

# NERSC Superfacility API



*Credits to:*

**LBLN Superfacility Team:**

D. Bard, C. Snavey, **G. Torok**, R. Thomas, A. Greiner, etc.

**NCEM:** P. Ercius, C. Harris (Kitware)

Bjoern Enders  
Data Science Workflows Architect  
NERSC/LBNL  
NERSC Data Day, Oct 26, 2022

# NERSC supports a large number of users and projects from DOE SC's experimental and observational facilities



Palomar Transient Factory  
Supernova



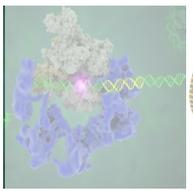
Planck Satellite  
Cosmic Microwave Background Radiation



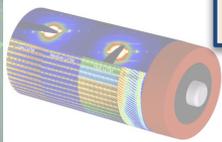
Dayabay Neutrinos



ALS Light Source



Cryo-EM



NCEM



DESI



LSST-DESC

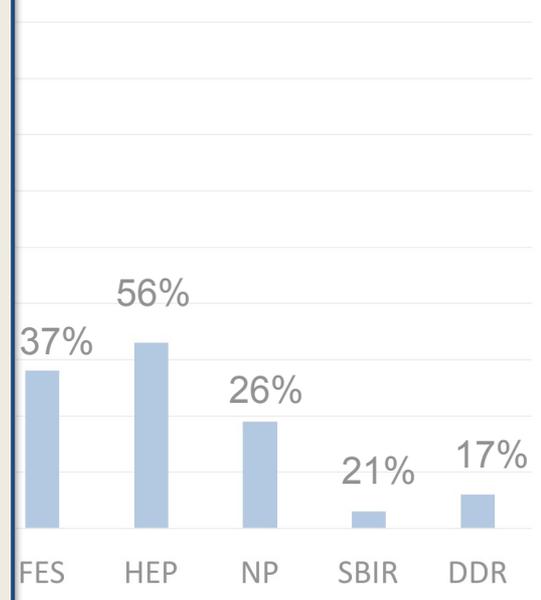


LZ

**Superfacility:**  
Ecosystem of connected facilities, software and expertise to enable new modes of discovery

**Superfacility API:**  
An API into NERSC to embed HPC into cross-facility workflows. It is also a general purpose API for all NERSC users and projects.

# of Projects Analyzing Experimental Data or Combining Modeling and Experimental Data by SC Office



~35% (235) of NERSC projects self identified as confirming the primary role of the project is to 1) analyze experimental data or; 2) create tools for experimental data analysis or; 3) combine experimental data with simulations and modeling

# Model case

Experiments at ext. facilities use high frame rate 2D detectors for their science.

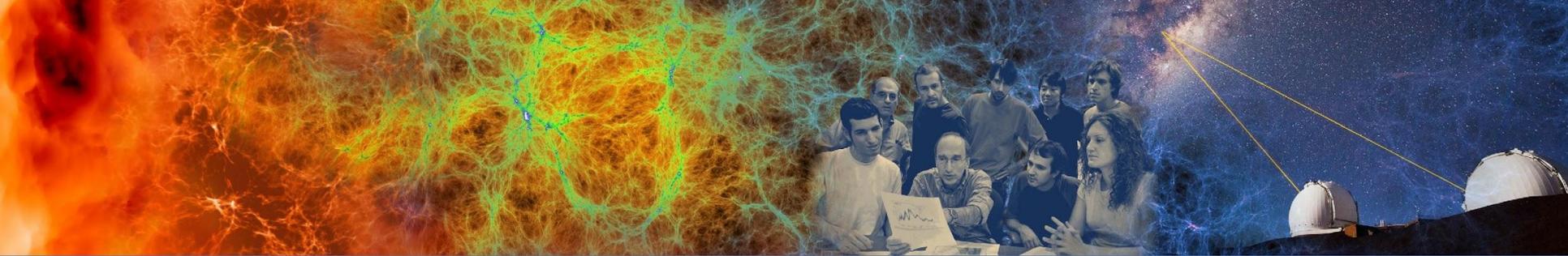
Hosting data & compute on site has become increasingly demanding.

## Requirements

- Planning (HPC as reliable partner)
  - machine-readable status
- Resiliency (needs failover)
  - compatible interfaces
- Realtime (can't wait in queue)
  - workflow endpoint
- Services (portals, data, db)



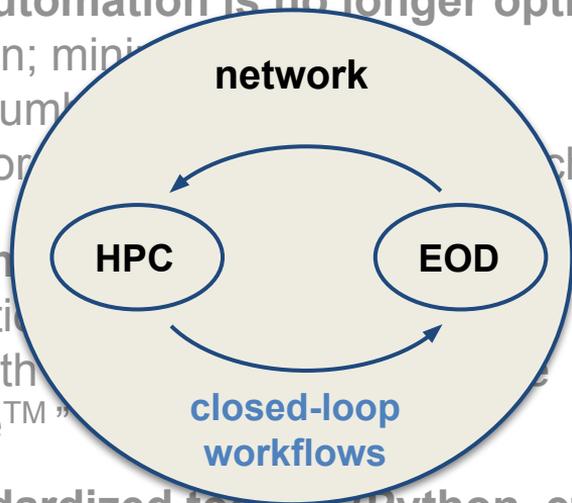
1. Plan / Check availability of NERSC resource for experiment.
  - check status / accounts
2. Get raw data to NERSC, when experiment is live.
  - move data
3. Start analysis job quasi synchronous with data
  - submit job / monitor job
4. Gather feedback, ideally immediate.
  - download / execute command
5. Move data and results to archive after analysis.
  - move data



# The API

# Why an API?

- **Meets a critical need; automation is no longer optional**
  - Unattended operation; minimize human intervention
  - Track/submit large number of jobs
  - Interface with collaborators and other machines
- **NERSC becomes “machine agnostic”**
  - Enables easier creation of workflows
  - Allows integration with other systems
    - “NERSC inside™”
- **Less DIY: simpler, standardized tooling (Python, etc)**
  - Stable refactor target for established projects or easier on-ramp for new ones
  - Contribute to HPC interface standards for portability
  - Authentication and security models



## Drivers:

- **Complex workflows**
- **Data-driven projects**
- **Real-time compute and streaming data from instruments**
- **Automation**

# What *specifically* can the API do?

**Vision: all NERSC interactions are callable;  
backend tools assist large or complex operations.**

## Endpoints prototyped or in prep:

- /status** query the status of NERSC component system health
- /account** data about the user's projects, roles, groups and usage information.
- /compute** run batch jobs, query job and queue statuses on compute resources.
- /storage** move data with Globus or between NERSC storage tiers
- /tasks** get info about asynchronous tasks (eg. from **/compute** or **/storage**).
- /utilities** traverse the filesystem, upload and download small files,  
and execute commands on NERSC systems
- /reservations** submit and manage future compute reservations (coming soon)

Action	Manual steps	With SuperFacility API
Check status	Test SSH or ping specific services for status	Query the <code>/status</code> API endpoint if resources are active.
Submit job	SSH in and submit jobs with <code>sbatch ...</code>	Create jobs using POST calls from a script or Spin service to the <code>/compute</code> endpoint.
Monitor job	SSH in (again) and do <code>queue   grep   sort   ...</code>	Consult the <code>/compute</code> and <code>/tasks</code> endpoints.
Plan ahead	Read the NERSC MOTD to see if any down time is planned	Query the <code>/status/outages/planned</code> API endpoint for planned outages
Move data	SSH in and run file transfer tools to move data	POST to the <code>/storage</code> API endpoint.
Check account	Log into "Iris" (our accounting web app) and check allocation account balance.	Query the <code>/account</code> API to get the same information.

# Model use

Experiments at ext. facilities use high frame rate 2D detectors for their science.

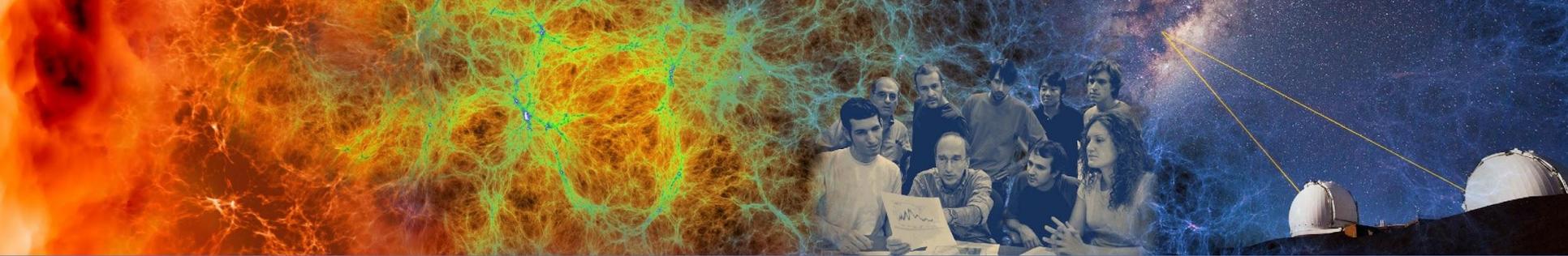
Hosting data & compute on site has become increasingly demanding.

## Requirements

- Planning (HPC as reliable partner)
  - machine-readable status
- Resiliency (needs failover)
  - compatible interfaces
- Realtime (can't wait in queue)
  - workflow endpoint
- Services (portals, data, db)



1. Plan / Check availability of NERSC resource for experiment.
  - `/status (/reservations)`
2. Get data to NERSC, when experiment is live.
  - `/storage`
3. Start analysis job quasi synchronous with data
  - `/compute /tasks`
4. Gather feedback, ideally immediate.
  - `/utilities /storage`
5. Move data and results to archive after analysis.
  - `/storage`



# A science example

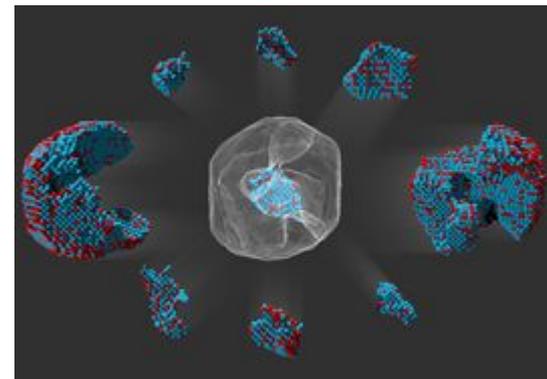


- The MF provides state-of-the-art expertise, methods, and instrumentation in nanoscale science in a safe environment free of charge
- NCEM is one of 7 facilities in the MF (about ~1/3 of total proposals and staff)
- Staff Scientists work in a 50/50 model: 50% of their time is spent on user research and 50% of their time is spent on internal research. User research is often highly collaborative.



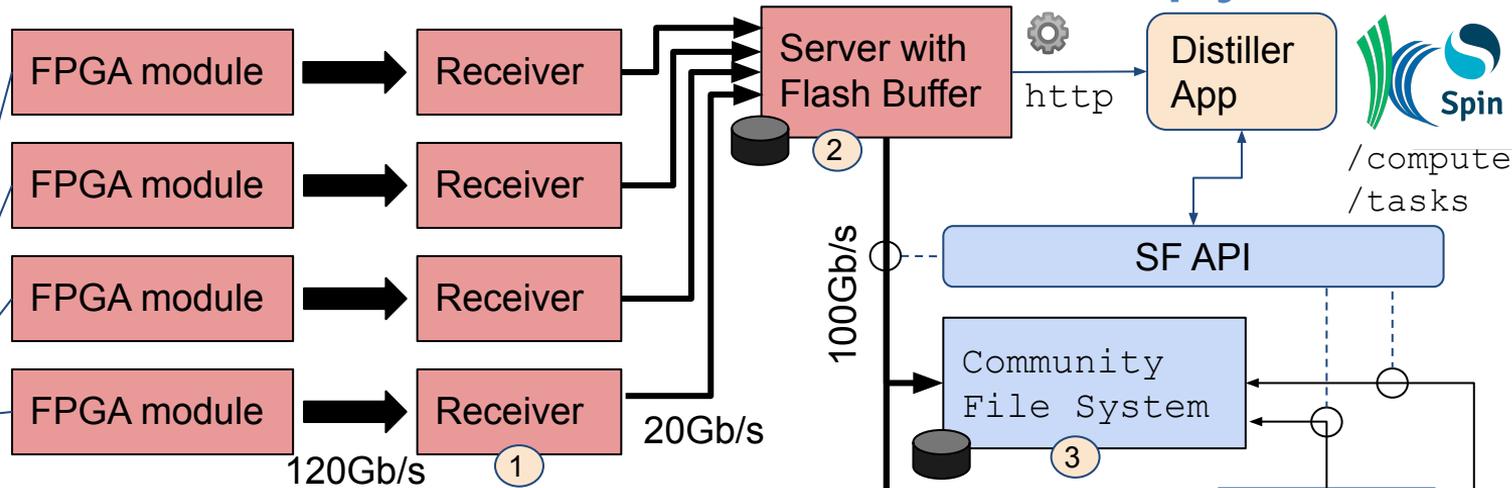
## We are leaders in:

- high resolution
- tomography
- in situ
- soft materials
- 4DSTEM
- image simulation
- electron detector technology



Y. Yang, C Ophus, M. Scott, P. Ercius, J. Miao, et al., *Nature* (2017)

# National Center for Electron Microscopy ...



... uses NERSC to process large data sets live during experiments

- 87,000 Hz (480 Gbit/s) readout (typical STEM scanning rates)
- 1kx1k scan is 650 GB captured in 15 seconds
- Data pipeline: FPGA → RAM → Flash storage → Sparse HDF5

15 sec (1) 140 sec (2) 5 min (3)

# NCEM - Distiller app

credits: Chris Harris @ Kitware, Peter Ercius @ NCEM

The screenshot shows the 'distiller.lbl.gov/scans/50' interface. On the left is a heatmap image. On the right, the scan details are displayed:

- Scan ID: 288
- Location: 128.55.132.192
- Created: 2021-10-20T12:38:39.093661
- Progress:
- Notes: MgFe single atom; 1kx1k

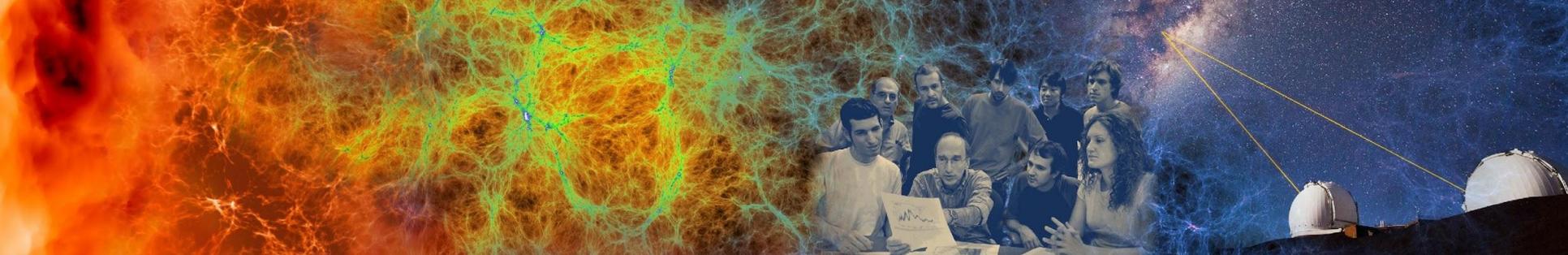
Buttons for 'TRANSFER' and 'COUNT' are visible. Below the scan details is a 'Jobs' table:

ID	Type	Slurm ID	State
49		48771314	<input checked="" type="checkbox"/>
42			

The screenshot shows the 'distiller.lbl.gov' interface with a list of scans. Each row includes a heatmap thumbnail, ID, Scan ID, Notes, Location, Created timestamp, and Progress status.

ID	Scan ID	Notes	Location	Created	Progress
50	288	MgFe single atom; 1kx1k	128.55.132.192	2021-10-20T12:38:39.093661	<input checked="" type="checkbox"/>
57	287	MgFe single atom	128.55.132.192	2021-10-20T12:20:40.864354	<input checked="" type="checkbox"/>
41	280	PiZn-SP	128.55.132.192	2021-10-19T15:41:48.427279	<input checked="" type="checkbox"/>
46	279	PiZn-SP 512x512	128.55.132.192	2021-10-19T15:09:01.756816	<input checked="" type="checkbox"/>
45	277	PiZn-SP	128.55.132.192	2021-10-19T14:17:45.569347	<input checked="" type="checkbox"/>
39	276	PiZn-SP	128.55.132.192	2021-10-19T14:09:30.814900	<input checked="" type="checkbox"/>
43	275	PiZn-SP	128.55.132.192	2021-10-19T14:01:24.432805	<input checked="" type="checkbox"/>
40	274	PiZn-SP	128.55.132.192	2021-10-19T13:56:50.663583	<input checked="" type="checkbox"/>
36	273	PiZn-SP	128.55.132.192	2021-10-19T13:50:25.172363	<input checked="" type="checkbox"/>
37	272	PiZn-SP	128.55.132.192	2021-10-19T13:48:21.832361	<input checked="" type="checkbox"/>
49	271	PiZn-SP	128.55.132.192	2021-10-19T13:01:54.504319	<input checked="" type="checkbox"/>
47	270	PiZn-SP	128.55.132.192	2021-10-19T13:01:07.076319	<input checked="" type="checkbox"/>
42	269	PiZn-SP	128.55.132.192	2021-10-19T13:00:45.756318	<input checked="" type="checkbox"/>

	ID	Scan ID	Notes	Location	Created	Progress
	81	14	STEMH some lattice	 128.55.132.192	2021-11-04T17:13:04.090986	<input checked="" type="checkbox"/>
	80	13	STEMH better alignment	 128.55.132.192	2021-11-04T17:10:50.558984	<input checked="" type="checkbox"/>
	79	12	STEMH 512x512 100p Au on edge	 128.55.132.192	2021-11-04T16:54:53.814969	<input checked="" type="checkbox"/>
	78	11	512x512 8mm CL	 128.55.132.192	2021-11-04T16:06:10.726925	<input checked="" type="checkbox"/>
	77	10	512x512 Au	 128.55.132.192	2021-11-04T15:54:21.766914	<input checked="" type="checkbox"/>
	76	9	512x512x4 Au	 128.55.132.192	2021-11-03T16:19:13.325631	<input checked="" type="checkbox"/>
	75	8	512x512x4 Au	 128.55.132.192	2021-11-03T15:51:44.517606	<input type="checkbox"/>
	74	7	128x128 Au	 128.55.132.192	2021-11-03T15:41:57.041597	<input checked="" type="checkbox"/>
	73	5		 128.55.132.192	2021-11-01T13:05:29.114843	<input checked="" type="checkbox"/>
	72	308		 128.55.132.192	2021-10-22T14:03:03.121835	<input checked="" type="checkbox"/>
	71	307		 128.55.132.192	2021-10-22T13:23:56.489800	<input checked="" type="checkbox"/>
	70	306		 128.55.132.192	2021-10-22T13:19:35.265796	<input checked="" type="checkbox"/>
	50	288	MgFe single atom; 1kx1k	 <span>cori</span> 128.55.132.192	2021-10-20T12:38:39.093661	<input checked="" type="checkbox"/>
	57	287	MgFe single atom	 128.55.132.192	2021-10-20T12:20:40.864354	<input checked="" type="checkbox"/>
	41	280	PtZn-SP	 <span>cori</span> 128.55.132.192	2021-10-19T15:41:48.427279	<input checked="" type="checkbox"/>



# How to use the API

# Superfacility API Basics

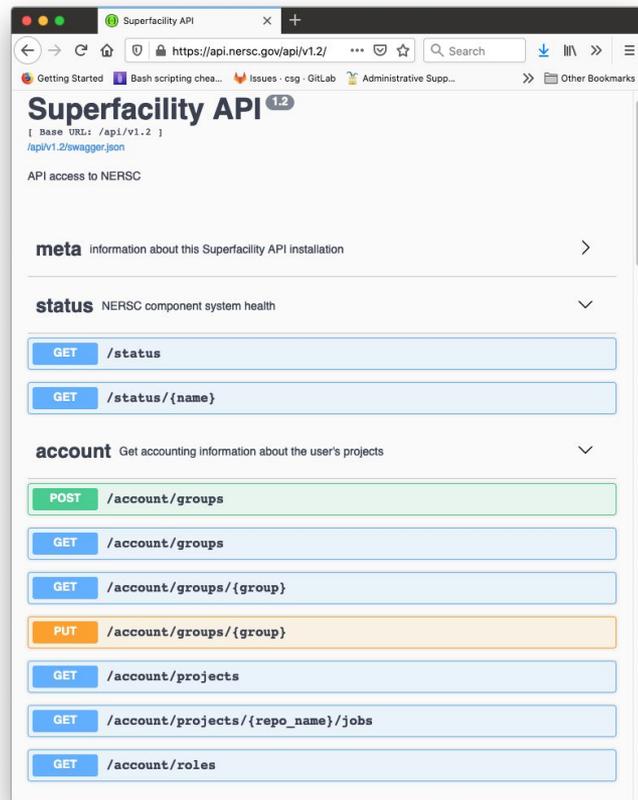
<https://api.nersc.gov/api/v1.2>

- A unified programmatic approach to accessing NERSC
- REST API with json input/output
- Standards-based authentication
- End user docs and examples:  
<https://docs.nersc.gov/services/sfapi/>

~ 4M logged requests since May,  
= one request every 3 sec

[Swagger](#) documentation:

- Interactive, up-to-date and self-documenting
- See endpoints, payloads, example code
- Works with any dev environment



# How to get your hands on the API

## As a user:

- The /status endpoints are all public.
- Every NERSC user can get API access to non-public endpoints via Iris <https://iris.nersc.gov>
  - Profile -> scroll down to "Superfacility API Clients" tab
  - R/W clients require filling out a form
- Getting started documentation available at <https://docs.nersc.gov/services/sfapi/>

## As an HPC facility:

- Please get in touch with us if you have question on how to build an API at your facility.
  - [benders@lbl.gov](mailto:benders@lbl.gov) (Bjoern Enders)
  - [djbard@lbl.gov](mailto:djbard@lbl.gov) (Debbie Bard)

### Create a New SuperFacility API Client ✕

**Workflow Type Client:**  
30 day lifetime (until 7/30/2021), read/write access  
This option is useful for long-running scripts such as a continuous automated workflow.

**Monitoring Type Client:**  
180 day lifetime (until 12/27/2021), read-only access  
This option is useful for when you only need read-only access, such as a monitoring dashboard.

**Client Name**

**Comments**

**Source IP Range (CIDR)**

IPv4 ranges must have /24 or greater suffix

# Example use

Experiments at ext. facilities use high frame rate 2D detectors for their science.

Hosting data & compute on site has become increasingly demanding.

## Requirements

- Planning (HPC as reliable partner)
  - [machine-readable status](#)
- Resiliency (needs failover)
  - [compatible interfaces](#)
- Realtime (can't wait in queue)
  - [workflow endpoint](#)
- Services (portals, data, db)



1. Plan / Check availability of NERSC resource for experiment.
  - [/status](#) ([/reservations](#))
2. Get data to NERSC, when experiment is live.
  - [/storage](#)
3. Start analysis job quasi synchronous with data
  - [/compute](#) [/tasks](#)
4. Gather feedback, ideally immediate.
  - [/utilities](#) [/storage](#)
5. Move data and results to archive after analysis.
  - [/storage](#)



```
Poll 04: {'id': '2236', 'status': 'completed', 'result': '{"status": "ok", "jobid": "484275", "error": null}'}
{'status': "ok", "jobid": "484275", "error": null}
```

Check queue status, wait for job to complete.

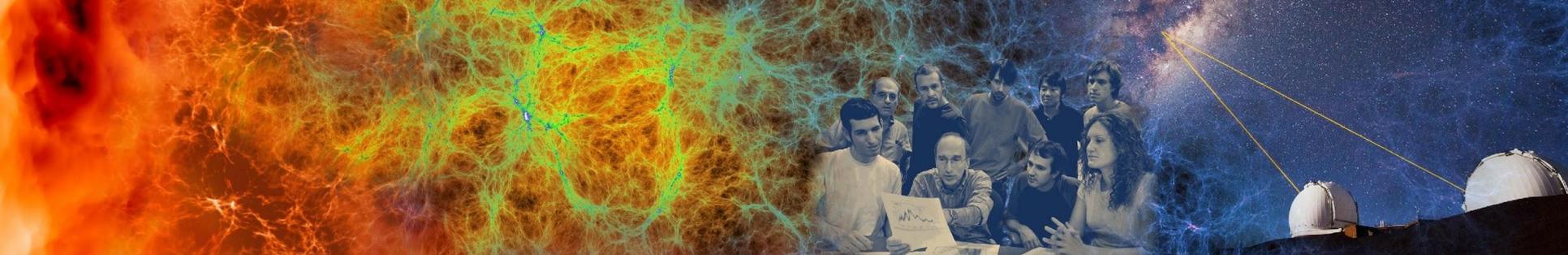
```
[338]: #print(json.dumps(api("compute/jobs/perlmutter/"+jobid+"?sacct=true"), indent=2))
print(api("compute/jobs/perlmutter/"+jobid+"?sacct=true")['output']['state'])
```

RUNNING

Read from the slurm output file

```
[334]: slurmfile = home+"/apidemo/slurm-"+jobid+".out"
response = api("utilities/command/perlmutter", { "executable": "tail -n 20 "+slurmfile })
if isinstance(response, Task):
    print(response.wait_for_result()['output'].strip())
```

```
Poll 03: {'id': '2237', 'status': 'completed', 'result': '{"status": "ok", "output": " * subpix          : linear\n\n * update_object_first : True\n\n * overlap_converge... : 0.05\n\n * overlap_max_itera... : 10\n\n * probe_inertia       : 1e-09\n\n * object_inertia      : 0.0001\n\n * fourier_power_bound : None\n\n * fourier_relax_factor : 0.05\n\n * obj_smooth_std      : None\n\n * clip_object         : None\n\n * probe_center_tol    : None\n\n * compute_log_likel... : True\n\n * probe_update_cuda... : False\n\n * object_update_cuda... : True\n\n * fft_lib             : reikna\n\n * alpha              : 1.0\n\n * name                : DM_pycuda\n\n=====\nIteration #10 of DM_pycuda :: Time 1.249\n\nErrors :: Fourier 5.94e+01, Photons 4.73e+01, Exit 3.51e+01\n", "error": null}'}
{'status': "ok", "output": " * subpix          : linear\n\n * update_object_first : True\n\n * overlap_converge... : 0.05\n\n * overlap_max_itera... : 10\n\n * probe_inertia       : 1e-09\n\n * object_inertia      : 0.0001\n\n * fourier_power_bound : None\n\n * fourier_relax_factor : 0.05\n\n * obj_smooth_std      : None\n\n * clip_object         : None\n\n * probe_center_tol    : None\n\n * compute_log_likel... : True\n\n * probe_update_cuda... : False\n\n * object_update_cuda... : True\n\n * fft_lib             : reikna\n\n * alpha              : 1.0\n\n * name                : DM_pycuda\n\n=====\nIteration #10 of DM_pycuda :: Time 1.249\n\nErrors :: Fourier 5.94e+01, Photons 4.73e+01, Exit 3.51e+01\n", "error": null}
* subpix          : linear
 * update_object_first : True
 * overlap_converge... : 0.05
 * overlap_max_itera... : 10
 * probe_inertia       : 1e-09
 * object_inertia      : 0.0001
 * fourier_power_bound : None
 * fourier_relax_factor : 0.05
 * obj_smooth_std      : None
 * clip_object         : None
 * probe_center_tol    : None
 * compute_log_likel... : True
```



# Roadmap

# Roadmap

- Clients and Tokens with more granular scope (~ weeks)
  - new client interface (draft see right)
  - more source IP ranges per client
  - short-lived full featured clients without manual review
- SF API to retire NEWT (~ months)
  - login-based route to get tokens for mynersc, science gateway apps or other web apps.
- Common API interface (~ year)
  - a set of endpoints/methods that work with many facilities
  - we're talking with CSCS (firecrest API), HPCS@LBL, S3DF@SLAC, OLCF

### Register a New Superfacility API Client

*Note: You don't need to register your client or use tokens if you only call endpoints that read API info or get system statuses.*

Client Name

Comments (optional)

User to create client for

Which security level does your client need?

Client credentials are scoped to enable endpoints by security level. Each level is valid for a certain length of time and number of IP address ranges. Choose the highest security level your application needs.

Green	Yellow	Orange	Red
60 days 16 IP ranges	60 days 8 IP ranges	30 days 8 IP ranges	2 days 2 IP ranges
Get user's projects Get user's account info Get user's roles Get user's filegroups Get info about a job Cancel a job List a directory Get status of a task Get statuses of tasks	All green functions + Get info about jobs Download a small file	All yellow functions + Create a group Get info about a group Update group members Start a transfer Upload a small file	All orange functions + Submit a job Run a command  Can be made valid for 30 days and 2 IP address ranges with <a href="#">security review</a> .

IP address range(s) (In CIDR format). Suffix must be /24 or higher.



# Outreach

- High-level overviews of the API have been given at workshops and meetings oriented toward software development, such as the DOE Workflow Workshop and Hack-a-thon and NERSC GPUs for Science Day, both in 2019.
- A proof-of-concept demonstrations with Jupyter notebooks were shown at the DOE exhibition booth at SC'19 and SC'21 (the latter already with Perlmutter)  
(<https://scdoe.info/demonstrations/>)
- A detailed presentation of the API architecture and usage coupled with a Jupyter-based demonstration was given at the Superfacility Project Demo Series in 2020.  
<https://www.youtube.com/watch?v=dmbBJmMUErU&list=PL20S5EeApOSsv6RVG6m0I6tx2wMp2T4PP&index=3>
- A [paper](#) was published with the ISC'21 SuperCompCloud workshop and accompanied by a presentation “[Automation for Data-Driven Research with the NERSC Superfacility API](#)”
- Science examples of earlier adopters were presented at a SC'21 BoF about HPC APIs.
  - Building an HPC API community.
- We're in touch with OLCF to adapt a similar API for their facility